

Stochastic Stability Analysis of Perturbed Learning Automata with Constant Step-Size in Strategic-Form Games^{*}

Georgios C. Chasparis

Software Competence Center Hagenberg GmbH, Softwarepark 21, A-4232 Hagenberg,
Austria
`georgios.chasparis@scch.at`,

Abstract. This paper considers a class of reinforcement-learning that belongs to the family of Learning Automata and provides a stochastic-stability analysis in strategic-form games. For this class of dynamics, convergence to pure Nash equilibria has been demonstrated only for the fine class of potential games. Prior work primarily provides convergence properties of the dynamics through stochastic approximations, where the asymptotic behavior can be associated with the limit points of an ordinary-differential equation (ODE). However, analyzing global convergence through the ODE-approximation requires the existence of a Lyapunov or a potential function, which naturally restricts the applicability of these algorithms to a fine class of games. To overcome these limitations, this paper introduces an alternative framework for analyzing stochastic-stability that is based upon an explicit characterization of the (unique) invariant probability measure of the induced Markov chain.

1 Introduction

Recently, multi-agent formulations have been utilized to tackle distributed optimization problems, since communication and computation complexity might be an issue in centralized optimization problems. In such formulations, decisions are usually taken in a repeated fashion, where agents select their next actions based on their *own* prior experience of the game.

The present paper discusses a class of reinforcement-learning dynamics, that belongs to the large family of Learning Automata [1,2], within the context of (non-cooperative) strategic-form games. In this class of dynamics, agents are repeatedly involved in a game with a fixed payoff-matrix, and they need to decide which action to play next having only access to their *own* prior actions

^{*} This work has been partially supported by the European Union grant EU H2020-ICT-2014-1 project RePhrase (No. 644235).

and payoffs. In Learning Automata, agents build their confidence over an action through repeated selection of this action and proportionally to the reward received from this action. Naturally, it has been utilized to analyze human-like (bounded) rationality [3].

Reinforcement learning has been applied in evolutionary economics, for modeling human and economic behavior [3,4,5,6,7]. It is also highly attractive to several engineering applications, since agents do not need to know neither the actions of the other agents, nor their own utility function. It has been utilized for system identification and pattern recognition [8], distributed network formation and coordination problems [9].

In strategic-form games, the main goal is to derive conditions under which convergence to Nash equilibria can be achieved. In social sciences, deriving such conditions may be important for justifying emergence of certain social phenomena. In engineering, convergence to Nash equilibria may also be desirable in distributed optimization problems, when the set of optimal solutions coincides with the set of Nash equilibria.

In Learning Automata, deriving conditions under which convergence to Nash equilibria is achieved may not be a trivial task. In particular, there are two main difficulties: a) excluding convergence to pure strategies that are *not* Nash equilibria, and b) excluding convergence to mixed strategy profiles. As it will be discussed in detail in a forthcoming Section 2, for some classes of (discrete-time) reinforcement-learning algorithms, convergence to non-Nash pure strategies may be achieved with positive probability. Moreover, excluding convergence to mixed strategy profiles may only be achieved under strong conditions in the utilities of the agents, (e.g., existence of a potential function).

In the present paper, we consider a class of (discrete-time) reinforcement-learning algorithms introduced in [9] that is closely related to existing algorithms for modeling human-like behavior, e.g., [3]. The main difference with prior reinforcement learning schemes lies in a) the step-size sequence, and b) the perturbation (or *mutations*) term. The step-size sequence is assumed constant, thus introducing a fading-memory effect of past experiences in each agent's strategy. On the other hand, the perturbation term introduces errors in the selection process of each agent. Both these two features can be used for designing a desirable asymptotic behavior.

We provide an analytical framework for deriving conclusions over the asymptotic behavior of the dynamics that is based on an explicit characterization of the invariant probability measure of the induced Markov chain. In particular, *we show that in all strategic-form games satisfying the Positive-Utility Property, the support of the invariant probability measure coincides with the set of pure strategy profiles.* This extends prior work where nonconvergence to mixed strategy profiles may only be excluded under strong conditions in the payoff matrix (e.g.,

existence of a potential function). A detailed discussion of the exact contributions of this paper is provided in the forthcoming Section 2. At the end of the paper, we also provide a brief discussion over how the proposed framework can be further utilized to provide a more detailed characterization of the stochastically stable states (e.g., excluding convergence to non-Nash pure strategy profiles). Due to space limitations, this analysis is not presented in this paper.

In the remainder of the paper, Section 2 presents a class of reinforcement-learning dynamics, related work and the main contribution of this paper. Section 3 provides the main result of this paper (Theorem 1), where the set of stochastically stable states is characterized. A short discussion is also provided over the significance of this result and how it can be utilized to provide further conclusions. Finally, Section 4 provides the technical derivation of the main result and Section 5 presents concluding remarks.

Notation:

- For a Euclidean topological space $\mathcal{X} \subset \mathbb{R}^n$, let $\mathcal{N}_\delta(x)$ denote the δ -neighborhood of $x \in \mathbb{R}^n$, i.e.,

$$\mathcal{N}_\delta(x) \doteq \{y \in \mathcal{X} : |x - y| < \delta\},$$

where $|\cdot|$ denotes the Euclidean distance.

- e_j denotes the *unit vector* in \mathbb{R}^n where its j th entry is equal to 1 and all other entries is equal to 0.
- $\Delta(n)$ denotes the *probability simplex* of dimension n , i.e.,

$$\Delta(n) \doteq \{x \in \mathbb{R}^n : x \geq 0, \mathbf{1}^\top x = 1\}.$$

- For some set A in a topological space \mathcal{Z} , let $\mathbb{I}_A : \mathcal{Z} \rightarrow \{0, 1\}$ denote the index function, i.e.,

$$\mathbb{I}_A(x) \doteq \begin{cases} 1 & \text{if } x \in A, \\ 0 & \text{else.} \end{cases}$$

- δ_x denotes the Dirac measure at x .
- Let A be a finite set and let any (finite) probability distribution $\sigma \in \Delta(|A|)$. The random selection of an element of A will be denoted $\text{rand}_\sigma[A]$. If $\sigma = (1/|A|, \dots, 1/|A|)$, i.e., it corresponds to the uniform distribution, the random selection will be denoted by $\text{rand}_{\text{unif}}[A]$.

2 Reinforcement Learning

2.1 Terminology

We consider the standard setup of finite strategic-form games. Consider a finite set of agents (or *players*) $\mathcal{I} = \{1, \dots, n\}$, and let each agent have a finite set of

actions \mathcal{A}_i . Let $\alpha_i \in \mathcal{A}_i$ denote any such action of agent i . The set of *action profiles* is the Cartesian product $\mathcal{A} \doteq \mathcal{A}_1 \times \cdots \times \mathcal{A}_n$ and let $\alpha = (\alpha_1, \dots, \alpha_n)$ be a representative element of this set. We will denote $-i$ to be the complementary set $\mathcal{I} \setminus i$ and often decompose an action profile as follows $\alpha = (\alpha_i, \alpha_{-i})$. The *payoff/utility function* of agent i is a mapping $u_i(\cdot) : \mathcal{A} \rightarrow \mathbb{R}$. A *strategic-form game* is defined by the triple $\langle \mathcal{I}, \mathcal{A}, \{u_i(\cdot)\}_i \rangle$.

For the remainder of the paper, we will be concerned with strategic-form games that satisfy the **Positive-Utility Property**.

Property 1 (Positive Utility Property). For any agent $i \in \mathcal{I}$ and any action profile $\alpha \in \mathcal{A}$, $u_i(\alpha) > 0$.

2.2 Reinforcement-learning algorithm

We consider a form of reinforcement learning that belongs to the general class of *learning automata* [2]. In learning automata, each agent updates a finite probability distribution $x_i \in \Delta(|\mathcal{A}_i|)$ representing its beliefs with respect to the most profitable action. The precise manner in which $x_i(t)$ changes at time t , depending on the performed action and the response of the environment, completely defines the reinforcement learning model.

The proposed reinforcement learning model is described in Table 1. At the first step, each agent i updates its action given its current strategy vector $x_i(t)$. Its selection is slightly perturbed by a perturbation (or *mutations*) factor $\lambda > 0$, such that, with a small probability λ agent i follows a uniform strategy (or, it *trembles*). At the second step, agent i evaluates its new selection by collecting a utility measurement, while in the last step, agent i updates its strategy vector given its new experience.

Here we identify actions \mathcal{A}_i with vertices of the simplex, $\{e_1, \dots, e_{|\mathcal{A}_i|}\}$. For example, if agent i selects its j th action at time t , then $e_{\alpha_i(t)} \equiv e_j$. Note that by letting the step-size ϵ to be sufficiently small and since the utility function $u_i(\cdot)$ is uniformly bounded in \mathcal{A} , $x_i(t) \in \Delta(|\mathcal{A}_i|)$ for all t .

In case $\lambda = 0$, the above update recursion will be referred to as the *unperturbed reinforcement learning*.

2.3 Related work

Erev-Roth type dynamics In prior reinforcement learning in games, analysis has been restricted to decreasing step-size sequences $\epsilon(t)$ and $\lambda = 0$. More specifically, in [3], the step-size sequence of agent i is $\epsilon_i(t) = 1/(ct^\nu + u_i(\alpha(t+1)))$ for some positive constant c and for $0 < \nu < 1$ (in the place of the constant step size ϵ of (2)). A comparative model is also used by [6], with $\epsilon_i(t) = 1/(V_i(t) + u_i(\alpha(t+1)))$, where $V_i(t)$ is the accumulated benefits of agent i up to

At fixed time instances $t = 1, 2, \dots$, and for each agent $i \in \mathcal{I}$, the following steps are executed recursively. Let $\alpha_i(t)$ and $x_i(t)$ denote the current action and strategy of agent i , respectively.

1. (**action update**) Agent $i \in \mathcal{I}$ selects a new action $\alpha_i(t+1)$ as follows:

$$\alpha_i(t+1) = \begin{cases} \text{rand}_{x_i(t)}[\mathcal{A}_i], & \text{with probability } 1 - \lambda, \\ \text{rnad}_{\text{unif}}[\mathcal{A}_i], & \text{with probability } \lambda, \end{cases} \quad (1)$$

for some small perturbation factor $\lambda > 0$.

2. (**evaluation**) Agent i applies its new action $\alpha_i(t+1)$ and retrieves a measurement of its utility function $u_i(\alpha(t+1)) > 0$.
3. (**strategy update**) Agent i revises its strategy vector $x_i \in \Delta(|\mathcal{A}_i|)$ as follows:

$$\begin{aligned} x_i(t+1) &= x_i(t) + \epsilon \cdot u_i(\alpha(t+1)) \cdot [e_{\alpha_i(t+1)} - x_i(t)] \\ &\doteq \mathcal{R}_i(\alpha(t+1), x_i(t)), \end{aligned} \quad (2)$$

for some constant step size $\epsilon > 0$.

Table 1. Perturbed Reinforcement Learning.

time t which gives rise to an urn process [5]. Some similarities are also shared with the Cross' learning model of [4], where $\epsilon(t) = 1$ and $u_i(\alpha(t)) \leq 1$, and its modification presented in [10], where $\epsilon(t)$, instead, is assumed decreasing.

The main difference of the proposed reinforcement-learning algorithm (Table 1) lies in the perturbation parameter $\lambda > 0$ which was first introduced and analyzed in [9]. A state-dependent perturbation term has also been investigated in [11]. The perturbation parameter may serve as an equilibrium selection mechanism, since *it excludes convergence to non-Nash action profiles*. It resolved one of the main issues of several (discrete-time) reinforcement-learning algorithms, that is the positive probability of convergence to non-Nash action profiles under some conditions in the payoff function and the step-size sequence.

This issue has also been raised by [12,6]. Reference [12] considered the model by [3] and showed that convergence to non-Nash pure strategy profiles can be excluded as long as $c > u_i(\alpha)$ for all $i \in \mathcal{I}$ and $\nu = 1$. On the other hand, convergence to non-Nash action profiles was not an issue with the urn model of [5] (as analyzed in [6]). However, the use of an urn-process type step-size sequence significantly reduces the applicability of the reinforcement learning scheme. In conclusion, the perturbation parameter $\lambda > 0$ may serve as a design tool for reinforcing convergence to Nash equilibria without necessarily employing an urn-process type step-size sequence. For engineering applications this is a desirable feature.

Although excluding convergence to non-Nash pure strategies can be guaranteed by using $\lambda > 0$, establishing convergence to pure Nash equilibria may still be an issue, since it further requires excluding convergence to mixed strategy profiles. As presented in [11], this can be guaranteed only under strong conditions in the payoff matrix. For example, as shown in [11, Proposition 8], excluding convergence to mixed strategy profiles requires a) the existence of a potential function, b) conditions over the second gradient of the potential function. Requiring the existence of a potential function considerably restricts the class of games where equilibrium selection can be described. Furthermore, condition (b) may not easily be verified in games of large number of players or actions.

Learning automata Certain forms of learning automata have been shown to converge to Nash equilibria in some classes of strategic-form games. For example, in [2], and for a generalized nonlinear reward-inaction scheme, convergence to Nash equilibrium strategies can be shown in identical interest games. Similar are the results presented in [13] for a linear reward-inaction scheme. These convergence results are restricted to games of payoffs in $[0, 1]$. Extension to a larger class of games is possible if *absolute monotonicity* (cf., [2, Definition 8.1]) is shown (similarly to the discussion in [11, Proposition 8]).

Reference [14] introduced a class of linear reward-inaction schemes in combination with a coordinated exploration phase so that convergence to the efficient Nash equilibrium is achieved. However, coordination of the exploration phase requires communication between the players.

Recently, work by the author [9] has introduced a new class of learning automata (namely, perturbed learning automata) which can be applied in games with no restriction in the payoff matrix. Furthermore, a small perturbation factor also influences the decisions of the players, through which convergence to non-Nash pure strategy profiles can be excluded. However, to demonstrate global convergence, a monotonicity condition still needs to be established [11].

Q-learning Similar questions of convergence to Nash equilibria also appear in alternative reinforcement learning formulations, such as approximate dynamic programming methodologies and *Q*-learning. However, this is usually accomplished under a stronger set of assumptions, which increases the computational complexity of the dynamics. For example, the Nash-Q learning algorithm of [15] addresses the problem of maximizing the discounted expected rewards for each agent by updating an approximation of the cost-to-go function (or *Q*-values). Alternative objectives may be used, such as the minimax criterion of [16]. However, it is indirectly assumed that agents need to have full access to the joint action space and the rewards received by the other agents.

More recently, reference [17] introduces a Q -learning scheme in combination with either adaptive play or better-reply dynamics in order to attain convergence to Nash equilibria in potential games [18] or weakly-acyclic games. However, this form of dynamics require that each player observes the actions selected by the other players, since a Q -value needs to be assigned in each joint action.

When the evaluation of the Q -values is totally independent, as in the individual Q -learning in [19], then convergence to Nash equilibria has been shown only for 2-player zero-sum games and 2-player partnership games with countably many Nash equilibria. Currently, there are no convergent results in games in multi-player games.

Payoff-based learning The aforementioned types of dynamics can be considered as a form of payoff-based learning dynamics, since adaptation is only governed by the perceived utility of the players. Recently, there have been several attempts for establishing convergence to Nash equilibria through alternative payoff-based learning dynamics, (see, e.g., the benchmark-based dynamics of [20], or the aspiration-based dynamics in [21]). For these type of dynamics, convergence to Nash equilibria can be established without requiring any strong monotonicity property (e.g., in multi-player weakly-acyclic games in [20]). However, an investigation is required with respect to the resulting convergence rates as compared to the dynamics incorporating policy iterations (e.g., the Erev-Roth type of dynamics or the learning automata discussed above).

2.4 Objective

This paper provides an analytical framework for analyzing convergence in multi-player strategic-form games when players implement a class of perturbed learning-automata. We wish to impose no strong monotonicity assumptions in the structure of the game (e.g., the existence of a potential function). We provide a characterization of the invariant probability measure of the induced Markov chain that shows that only the pure-strategy profiles belong to its support. Thus, we implicitly exclude convergence to any mixed strategy profile (including mixed Nash equilibria). This result imposes no restrictions in the payoff matrix other than the Positive-Utility Property.

3 Convergence Analysis

3.1 Terminology and notation

Let $\mathcal{Z} \doteq \mathcal{A} \times \mathbf{\Delta}$, where $\mathbf{\Delta} \doteq \Delta(|\mathcal{A}_1|) \times \dots \times \Delta(|\mathcal{A}_n|)$, i.e., pairs of joint actions α and nominal strategy profiles x . The set \mathcal{A} is endowed with the discrete topology,

Δ with its usual Euclidean topology, and \mathcal{Z} with the corresponding product topology. We also let $\mathfrak{B}(\mathcal{Z})$ denote the Borel σ -field of \mathcal{Z} , and $\mathfrak{P}(\mathcal{Z})$ the set of probability measures on $\mathfrak{B}(\mathcal{Z})$ endowed with the Prohorov topology, i.e., the topology of weak convergence. The algorithm introduced in Table 1 defines an \mathcal{Z} -valued Markov chain. Let $P_\lambda : \mathcal{Z} \times \mathfrak{B}(\mathcal{Z}) \rightarrow [0, 1]$ denote its transition probability function (t.p.f.), parameterized by $\lambda > 0$. We refer to the process with $\lambda > 0$ as the *perturbed process*. Let also $P : \mathcal{Z} \times \mathfrak{B}(\mathcal{Z})$ denote the t.p.f. of the *unperturbed process*, i.e., when $\lambda = 0$.

We let $C_b(\mathcal{Z})$ denote the Banach space of real-valued continuous functions on \mathcal{Z} under the sup-norm (denoted by $\|\cdot\|_\infty$) topology. For $f \in C_b(\mathcal{Z})$, define

$$P_\lambda f(z) \doteq \int_{\mathcal{Z}} P_\lambda(z, dy) f(y),$$

and

$$\mu[f] \doteq \int_{\mathcal{Z}} \mu(dx) f(x), \text{ for } \mu \in \mathfrak{P}(\mathcal{Z}).$$

The process governed by the unperturbed process P will be denoted by $\{Z_t : t \geq 0\}$. Let $\Omega \doteq \mathcal{Z}^\infty$ denote the canonical path space, i.e., an element $\omega \in \Omega$ is a sequence $\{\omega(0), \omega(1), \dots\}$, with $\omega(t) = (\alpha(t), x(t)) \in \mathcal{Z}$. We use the same notation for the elements (α, x) of the space \mathcal{Z} and for the coordinates of the process $Z_t = (\alpha(t), x(t))$. Let also $\mathbb{P}_z[\cdot]$ denote the unique probability measure induced by the unperturbed process P on the product σ -algebra of \mathcal{Z}^∞ , initialized at $z = (\alpha, x)$, and $\mathbb{E}_z[\cdot]$ the corresponding expectation operator. Let also \mathfrak{F}_t , $t \geq 0$, denote the σ -algebra generated by $\{Z_\tau, \tau \leq t\}$.

3.2 Stochastic stability

First, we note that both P and P_λ ($\lambda > 0$) satisfy the *weak Feller property* (cf., [22, Definition 4.4.2]).

Proposition 1. *Both the unperturbed process P ($\lambda = 0$) and the perturbed process P_λ ($\lambda > 0$) have the weak Feller property.*

Proof. Let us consider any sequence $\{Z^{(k)} = (\alpha^{(k)}, x^{(k)})\}$ such that $Z^{(k)} \rightarrow Z = (\alpha, x) \in \mathcal{Z}$.

For the unperturbed process governed by $P(\cdot, \cdot)$, and for any open set $O \in \mathfrak{B}(\mathcal{Z})$, the following holds:

$$\begin{aligned} P(Z^{(k)} = (\alpha^{(k)}, x^{(k)}), O) \\ &= \sum_{\alpha \in \mathcal{P}_{\mathcal{A}}(O)} \left\{ \mathbb{P}_{Z^{(k)}}[\text{rand}_{x_i^{(k)}}[\mathcal{A}_i] = \alpha_i, \forall i \in \mathcal{I}] \cdot \right. \\ &\quad \left. \prod_{i=1}^n \mathbb{P}_{Z^{(k)}}[\mathcal{R}_i(\alpha, x_i^{(k)}) \in \mathcal{P}_{\mathcal{X}_i}(O)] \right\} \end{aligned}$$

$$= \sum_{\alpha \in \mathcal{P}_{\mathcal{A}}(O)} \left\{ \prod_{i=1}^n \mathbb{I}_{\mathcal{P}_{\mathcal{X}_i}(O)}(\mathcal{R}_i(\alpha, x_i^{(k)})) x_{i\alpha_i}^{(k)} \right\},$$

where $\mathcal{P}_{\mathcal{X}_i}(O)$ and $\mathcal{P}_{\mathcal{A}}(O)$ are the *canonical projections* defined by the product topology. Similarly, we have:

$$\begin{aligned} P(Z = (\alpha, x), O) \\ = \sum_{\alpha \in \mathcal{P}_{\mathcal{A}}(O)} \left\{ \prod_{i=1}^n \mathbb{I}_{\mathcal{P}_{\mathcal{X}_i}(O)}(\mathcal{R}_i(\alpha, x_i)) x_{i\alpha_i} \right\}. \end{aligned}$$

(a) Consider the case $x \in \Delta^\circ$, i.e., x belongs to the interior of Δ . For all $i \in \mathcal{I}$, due to the continuity of $\mathcal{R}_i(\cdot, \cdot)$ with respect to its second argument, and the fact that O is an open set, there exists $\delta > 0$ such that $\mathbb{I}_{\mathcal{P}_{\mathcal{X}_i}(O)}(\mathcal{R}_i(\alpha, x_i)) = \mathbb{I}_{\mathcal{P}_{\mathcal{X}_i}(O)}(\mathcal{R}_i(\alpha, y_i))$ for all $y_i \in \mathcal{N}_\delta(x_i)$. Thus, for any sequence $Z^{(k)} = (\alpha^{(k)}, x^{(k)})$ such that $Z^{(k)} \rightarrow Z = (\alpha, x)$, we have that $P(Z^{(k)}, O) \rightarrow P(Z, O)$, as $k \rightarrow \infty$.

(b) Consider the case $x \in \partial\Delta$, i.e., x belongs to the boundary of Δ . Then, there exists $i \in \mathcal{I}$ such that $x_i \in \partial\Delta(|\mathcal{A}_i|)$, i.e., there exists an action $j \in \mathcal{A}_i$ such that $x_{ij} = 0$. For any open set $O \in \mathfrak{B}(\mathcal{Z})$, $x_i \notin \mathcal{P}_{\mathcal{X}_i}(O)$. Furthermore, for any $\alpha_i \in \text{rand}_{x_i}[\mathcal{A}_i]$, $\mathbb{I}_{\mathcal{P}_{\mathcal{X}_i}(O)}(\mathcal{R}_i((\alpha_i, \alpha_{-i}), x_i)) = 0$ (since $x_{ij} = 0$ and therefore x_i cannot escape from the boundary). This directly implies that $P(Z = (\alpha, x), O) = 0$. Construct a sequence $(\alpha^{(k)}, x^{(k)})$ that converges to (α, x) such that $\alpha^{(k)} = \alpha$, $x_{i\alpha_i}^{(k)} > 0$ and $x_i = e_{\alpha_i}$, i.e., the strategy of player i converges to the vertex of action α_i . Pick also $O \in \mathfrak{B}(\mathcal{Z})$, such that $\mathbb{I}_{\mathcal{P}_{\mathcal{X}_i}(O)}(\mathcal{R}_i(\alpha, x_i^{(k)})) = 1$ for all large k . This is always possible by selecting an open set O such that $x \in \partial\mathcal{P}_{\mathcal{X}}(O)$ and $x^{(k)} \in \mathcal{P}_{\mathcal{X}}(O)$ for all k . In this case, $\lim_{k \rightarrow \infty} P(Z^{(k)}, O) = 1$. We conclude that for any sequence $Z^{(k)} = (\alpha^{(k)}, x^{(k)})$ that converges to $Z = (\alpha, x)$, such that $x \in \partial\Delta$, and for any open set $O \in \mathfrak{B}(\mathcal{Z})$,

$$\lim_{k \rightarrow \infty} P(Z^{(k)}, O) \geq P(Z, O) = 0.$$

By [22, Proposition 7.2.1], we conclude that P satisfies the weak Feller property. The same steps can be followed to show that P_λ also satisfies the weak Feller property. •

The measure $\mu_\lambda \in \mathfrak{P}(\mathcal{Z})$ is called an *invariant probability measure* for P_λ if

$$(\mu_\lambda P_\lambda)(A) \doteq \int_{\mathcal{Z}} \mu_\lambda(dx) P_\lambda(z, A) = \mu_\lambda(A), \quad A \in \mathfrak{B}(\mathcal{Z}).$$

Since \mathcal{Z} defines a locally compact separable metric space and P, P_λ have the weak Feller property, they both admit an invariant probability measure, denoted μ and μ_λ , respectively [22, Theorem 7.2.3].

We would like to characterize the *stochastically stable states* $z \in \mathcal{Z}$ of P_λ , that is any state $z \in \mathcal{Z}$ for which any collection of invariant probability measures $\{\mu_\lambda \in \mathfrak{P}(\mathcal{Z}) : \mu_\lambda P_\lambda = \mu_\lambda, \lambda > 0\}$ satisfies $\liminf_{\lambda \rightarrow 0} \mu_\lambda(z) > 0$. As the

forthcoming analysis will show, the stochastically stable states will be a subset of the set of *pure strategy states* (p.s.s.) defined as follows:

Definition 1 (Pure Strategy State). *A pure strategy state is a state $s = (\alpha, x) \in \mathcal{Z}$ such that for all $i \in \mathcal{I}$, $x_i = e_{\alpha_i}$, i.e., x_i coincides with the vertex of the probability simplex $\Delta(|\mathcal{A}_i|)$ which assigns probability 1 to action α_i .*

We will denote the set of pure strategy states by \mathcal{S} .

Theorem 1 (Stochastic Stability). *There exists a unique probability vector $\pi = (\pi_1, \dots, \pi_{|\mathcal{S}|})$ such that for any collection of invariant probability measures $\{\mu_\lambda \in \mathfrak{P}(\mathcal{Z}) : \mu_\lambda P_\lambda = \mu_\lambda, \lambda > 0\}$, the following hold:*

- (a) $\lim_{\lambda \rightarrow 0} \mu_\lambda(\cdot) = \hat{\mu}(\cdot) \doteq \sum_{s \in \mathcal{S}} \pi_s \delta_s(\cdot)$, where convergence is in the weak sense.
- (b) The probability vector π is an invariant distribution of the (finite-state) Markov process \hat{P} , such that, for any $s, s' \in \mathcal{S}$,

$$\hat{P}_{ss'} \doteq \lim_{t \rightarrow \infty} QP^t(s, \mathcal{N}_\delta(s')), \quad (3)$$

for any $\delta > 0$ sufficiently small, where Q is the t.p.f. corresponding to only one player trembling (i.e., following the uniform distribution of (1)).

The proof of Theorem 1 requires a series of propositions and will be presented in detail in Section 4.

3.3 Discussion

Theorem 1 establishes an important observation. That is, the “*equivalence*” (in a weak convergence sense) of the original (perturbed) learning process with a simplified process, where *agents simultaneously tremble at the first iteration and then they do not tremble*. This form of simplification of the dynamics has originally been exploited to analyze *aspiration learning* dynamics in [21], and it is based upon the fact that *under the unperturbed dynamics, agents’ strategies will eventually converge to a pure strategy profile*.

Furthermore, the limiting behavior of the original (perturbed) dynamics can be characterized by the (*unique*) invariant distribution of a finite-state Markov chain $\{P_{ss'}\}$, whose states correspond to the pure-strategy states of the game. In other words, *we should expect that as the perturbation parameter λ approaches zero, the algorithm spends the majority of the time on pure strategy profiles*. The importance of this result lies on the fact that no constraints have been imposed in the payoff matrix of the game other than the Positive-Utility Property 1. Thus, it extends to games beyond the fine set of potential games.

This convergence result can further be augmented with an ODE analysis for stochastic approximations to exclude convergence to pure strategies that are not Nash equilibria (as derived in [11] for the case of diminishing step size). Due to space limitations this analysis is not presented in this paper, however it can be the subject of future work.

4 Technical Derivation

4.1 Unperturbed Process

For $t \geq 0$ define the sets

$$A_t \doteq \{\omega \in \Omega : \alpha(\tau) = \alpha(t), \text{ for all } \tau \geq t\},$$

$$B_t \doteq \{\omega \in \Omega : \alpha(\tau) = \alpha(0), \text{ for all } 0 \leq \tau \leq t\}.$$

Note that $\{B_t : t \geq 0\}$ is a non-increasing sequence, i.e., $B_{t+1} \subseteq B_t$, while $\{A_t : t \geq 0\}$ is non-decreasing, i.e., $A_{t+1} \supseteq A_t$. Let

$$A_\infty \doteq \bigcup_{t=0}^{\infty} A_t \text{ and } B_\infty \doteq \bigcap_{t=1}^{\infty} B_t.$$

In other words, *the set A_∞ corresponds to the event that agents eventually play the same action profile, while B_∞ corresponds to the event that agents never change their actions.*

Proposition 2 (Convergence to p.s.s.). *Let us assume that the step size $\epsilon > 0$ is sufficiently small such that $0 < \epsilon u_i(\alpha) < 1$ for all $\alpha \in \mathcal{A}$ and for all agents $i \in \mathcal{I}$. Then, the following hold:*

- (a) $\inf_{z \in \mathcal{Z}} \mathbb{P}_z[B_\infty] > 0$,
- (b) $\inf_{z \in \mathcal{Z}} \mathbb{P}_z[A_\infty] = 1$.

The first statement of Proposition 2 states that *the probability that agents never change their actions is bounded away from zero*, while the second statement states that *the probability that eventually agents play the same action profile is one*.

Proof. (a) Let us consider an action profile $\alpha = (\alpha_1, \dots, \alpha_n) \in \mathcal{A}$, and an initial strategy profile $x(0) = (x_1(0), \dots, x_n(0))$ such that $x_{i\alpha_i}(0) > 0$ for all $i \in \mathcal{I}$. Note that if the same action profile α is selected up to time t , then the strategy of agent i satisfies:

$$x_i(t) = e_{\alpha_i} - (1 - \epsilon u_i(\alpha))^t (e_{\alpha_i} - x_i(0)). \quad (4)$$

Given that B_t is non-increasing, from continuity from above we have

$$\mathbb{P}_z[B_\infty] = \lim_{t \rightarrow \infty} \mathbb{P}_z[B_t] = \lim_{t \rightarrow \infty} \prod_{k=0}^t \prod_{i=1}^n x_{i\alpha_i}(k).$$

Note that $\mathbb{P}[B_\infty] > 0$ if and only if

$$\sum_{t=1}^{\infty} \log(x_{i\alpha_i}(t)) > -\infty. \quad (5)$$

Let us introduce the variable

$$y_i(t) \doteq 1 - x_{i\alpha_i}(t) = \sum_{j \in \mathcal{A}_i \setminus \alpha_i} x_{ij}(t),$$

which corresponds to the probability of agent i selecting any action other than α_i . Condition (5) is equivalent to

$$-\sum_{t=0}^{\infty} \log(1 - y_i(t)) < \infty, \text{ for all } i \in \mathcal{I}. \quad (6)$$

We also have that

$$\lim_{t \rightarrow \infty} \frac{-\log(1 - y_i(t))}{y_i(t)} = \lim_{t \rightarrow \infty} \frac{1}{1 - y_i(t)} > \rho$$

for some $\rho > 0$, since $0 \leq y_i(t) \leq 1$. Thus, from the Limit Comparison Test, we conclude that condition (6) holds if and only if $\sum_{t=1}^{\infty} y_i(t) < \infty$, for each $i \in \mathcal{I}$.

Lastly, note that $y_i(t+1)/y_i(t) = 1 - \epsilon u_i(\alpha)$. By Raabe's criterion, the series $\sum_{t=0}^{\infty} y_i(t)$ is convergent if $\lim_{t \rightarrow \infty} t(y_i(t)/y_i(t+1) - 1) > 1$. We have

$$t \left(\frac{y_i(t)}{y_i(t+1)} - 1 \right) = t \frac{\epsilon u_i(\alpha)}{1 - \epsilon u_i(\alpha)}.$$

Thus, if $\epsilon u_i(\alpha) < 1$ for all $\alpha \in \mathcal{A}$ and $i \in \mathcal{I}$, then $1 - \epsilon u_i(\alpha) > 0$ and $\lim_{t \rightarrow \infty} t(\epsilon u_i(\alpha)/(1 - \epsilon u_i(\alpha))) > 1$, which implies that the series $\sum_{t=1}^{\infty} y_i(t)$ is convergent. Thus, we conclude that $\mathbb{P}_z[B_{\infty}] > 0$.

(b) Define the event

$$C_t \doteq \{ \exists \alpha' \neq \alpha(t) : x_{i\alpha'_i}(t) > 0, \text{ for all } i \in \mathcal{I} \},$$

i.e., C_t corresponds to the event that there exists an action profile different from the current action profile for which the nominal strategy assigns positive probability for all agents i . Note that $A_t^c \subseteq C_t$, since A_t^c occurs only if there is some action profile $\alpha' \neq \alpha(t)$ for which the nominal strategy assigns positive probability. This further implies that $\mathbb{P}_z[A_t^c] \leq \mathbb{P}_z[C_t]$. Then, we have:

$$\begin{aligned} & \mathbb{P}_z[A_{t+1}|A_t^c] \\ &= \frac{\mathbb{P}_z[A_{t+1} \cap A_t^c]}{\mathbb{P}_z[A_t^c]} \\ &\geq \frac{\mathbb{P}_z[A_{t+1} \cap A_t^c]}{\mathbb{P}_z[C_t]} \\ &\geq \mathbb{P}_z[A_{t+1} \cap A_t^c | C_t] \\ &= \mathbb{P}_z[\{\alpha(\tau) = \alpha' \neq \alpha(t), \forall \tau > t\} | C_t] \\ &\geq \inf_{\alpha' \neq \alpha} \prod_{i=1}^n x_{i\alpha'_i}(t) \prod_{k=t+1}^{\infty} \{1 - (1 - \epsilon u_i(\alpha'))^{k-t-1} c_i(\alpha')\} \end{aligned}$$

$$\geq \inf_{\alpha' \neq \alpha} \prod_{i=1}^n x_{i\alpha'_i}(t) \prod_{k=0}^{\infty} \{1 - (1 - \epsilon u_i(\alpha'))^k c_i(\alpha')\}$$

where $c_i(\alpha') \doteq 1 - x_{i\alpha'_i}(t) \geq 0$. We have already shown in part (a) that the second part of the r.h.s. is bounded away from zero. Therefore, we conclude that $\mathbb{P}_z[A_{t+1}|A_t^c] > 0$. Thus, from the counterpart of the Borel-Cantelli Lemma, $\mathbb{P}_z[A_\infty] = 1$. •

The above proposition is rather useful in characterizing the support of any invariant measure of the unperturbed process, as the following proposition shows.

Proposition 3 (Limiting t.p.f. of unperturbed process). *Let μ denote an invariant probability measure of P . Then, there exists a t.p.f. Π on $\mathcal{Z} \times \mathfrak{B}(\mathcal{Z})$ such that*

- (a) *for μ -a.e. $z \in \mathcal{Z}$, $\Pi(z, \cdot)$ is an invariant probability measure for P ;*
- (b) *for all $f \in C_b(\mathcal{Z})$, $\lim_{t \rightarrow \infty} \|P^t f - \Pi f\|_\infty = 0$;*
- (c) *μ is an invariant probability measure of Π ;*
- (d) *the support¹ of Π is on \mathcal{S} for all $z \in \mathcal{Z}$.*

Proof. The state space \mathcal{Z} is a locally compact separable metric space and the t.p.f. of the unperturbed process P admits an invariant probability measure due to Proposition 1. Thus, statements (a), (b) and (c) follow directly from [22, Theorem 5.2.2 (a), (b), (e)].

(d) Let us assume that the support of Π includes points in \mathcal{Z} other than the pure strategy states. Let also $O \subset \mathcal{Z}$ be an open set such that $O \cap \mathcal{S} = \emptyset$ and $\Pi(z^*, O) > 0$ for some $z^* \in \mathcal{Z}$. Given that P^t converges weakly to Π as $t \rightarrow \infty$, from Portmanteau theorem (cf., [22, Theorem 1.4.16]), we have that

$$\liminf_{t \rightarrow \infty} P^t(z^*, O) \geq \Pi(z^*, O) > 0.$$

This is a contradiction of Proposition 2(b). Thus, the conclusion follows. •

Proposition 3 states that the limiting unperturbed t.p.f. converges weakly to a t.p.f. Π which accepts the same invariant p.m. as P . Furthermore, *the support of Π is the set of pure strategy states \mathcal{S}* . This is a rather important observation, since the limiting perturbed process can also be “related” (in a weak-convergence sense) to the t.p.f. Π , as it will be shown in the following section.

¹ The *support* of a measure μ on \mathcal{Z} is the unique closed set $F \subset \mathfrak{B}(\mathcal{Z})$ such that $\mu(\mathcal{Z} \setminus F) = 0$ and $\mu(F \cap O) > 0$ for every open set $O \subset \mathcal{Z}$ such that $F \cap O \neq \emptyset$.

4.2 Decomposition of perturbed t.p.f.

We can decompose the t.p.f. of the perturbed process as follows:

$$P_\lambda = (1 - \varphi(\lambda))P + \varphi(\lambda)Q_\lambda$$

where $\varphi(\lambda) = 1 - (1 - \lambda)^n$ is the probability that at least one agent trembles (since $(1 - \lambda)^n$ is the probability that no agent trembles), and Q_λ corresponds to the t.p.f. induced by the one-step reinforcement-learning update when at least one agent trembles. Note that $\varphi(\lambda) \rightarrow 0$ as $\lambda \rightarrow 0$.

Define also Q to be the t.p.f. when *only one* players trembles, and Q^* is the t.p.f. where at least two players tremble. Then, we may write:

$$Q_\lambda = (1 - \psi(\lambda))Q + \psi(\lambda)Q^*, \quad (7)$$

where $\psi(\lambda) \doteq 1 - \frac{n\lambda}{1-(1-\lambda)^n}$ corresponds to the probability that at least two players tremble given that at least one player trembles.

Let us also define the infinite-step t.p.f. when trembling only at the first step (briefly, *lifted* t.p.f.) as follows:

$$P_\lambda^L \doteq \varphi(\lambda) \sum_{t=0}^{\infty} (1 - \varphi(\lambda))^t Q_\lambda P^t = Q_\lambda R_\lambda \quad (8)$$

where $R_\lambda \doteq \varphi(\lambda) \sum_{t=0}^{\infty} (1 - \varphi(\lambda))^t P^t$, i.e., R_λ corresponds to the *resolvent* t.p.f.

Proposition 4 (Invariant p.m. of perturbed process). *The following hold:*

- (a) For $f \in C_b(\mathcal{Z})$, $\lim_{\lambda \rightarrow 0} \|R_\lambda f - \Pi f\|_\infty = 0$.
- (b) For $f \in C_b(\mathcal{Z})$, $\lim_{\lambda \rightarrow 0} \|P_\lambda^L f - Q \Pi f\|_\infty = 0$.
- (c) Any invariant distribution μ_λ of P_λ is also an invariant distribution of P_λ^L .
- (d) Any weak limit point in $\mathfrak{P}(\mathcal{Z})$ of μ_λ , as $\lambda \rightarrow 0$, is an invariant probability measure of $Q \Pi$.

Proof. (a) For any $f \in C_b(\mathcal{Z})$, we have

$$\begin{aligned} & \|R_\lambda f - \Pi f\|_\infty \\ &= \left\| \varphi(\lambda) \sum_{t=0}^{\infty} (1 - \varphi(\lambda))^t P^t f - \Pi f \right\|_\infty \\ &= \left\| \varphi(\lambda) \sum_{t=0}^{\infty} (1 - \varphi(\lambda))^t (P^t f - \Pi f) \right\|_\infty \end{aligned}$$

where we have used the property $\varphi(\lambda) \sum_{t=0}^{\infty} (1 - \varphi(\lambda))^t = 1$. Note that

$$\varphi(\lambda) \sum_{t=0}^{\infty} (1 - \varphi(\lambda))^t \|P^t f - \Pi f\|_\infty$$

$$\leq (1 - \varphi(\lambda))^T \sup_{t \geq T} \|P^t f - \Pi f\|_\infty.$$

From Proposition 3(b), we have that for any $\delta > 0$, there exists $T = T(\delta) > 0$ such that the r.h.s. is uniformly bounded by δ for all $t \geq T$. Thus, the sequence

$$A_T \doteq \varphi(\lambda) \sum_{t=0}^T (1 - \varphi(\lambda))^t (P^t f - \Pi f)$$

is Cauchy and therefore convergent (under the sup-norm). In other words, there exists $A \in \mathbb{R}$ such that

$$\lim_{T \rightarrow \infty} \|A_T - A\|_\infty = 0.$$

For every $T > 0$, we have

$$\|R_\lambda f - \Pi f\|_\infty \leq \|A_T\|_\infty + \|A - A_T\|_\infty.$$

Note that

$$\|A_T\|_\infty \leq \varphi(\lambda) \sum_{t=0}^T (1 - \varphi(\lambda))^t \|P^t f - \Pi f\|_\infty.$$

If we take $\lambda \downarrow 0$, then the r.h.s. converges to zero. Thus,

$$\|R_\lambda f - \Pi f\|_\infty \leq \|A - A_T\|_\infty, \text{ for all } T > 0,$$

which concludes the proof.

(b) For any $f \in C_b(\mathcal{Z})$, we have

$$\begin{aligned} & \|P_\lambda^L f - Q \Pi f\|_\infty \\ & \leq \|Q_\lambda (R_\lambda f - \Pi f)\|_\infty + \|Q_\lambda \Pi f - Q \Pi f\|_\infty \\ & \leq \|R_\lambda f - \Pi f\|_\infty + \|Q_\lambda \Pi f - Q \Pi f\|_\infty. \end{aligned}$$

The first term of the r.h.s. approaches 0 as $\lambda \downarrow 0$ according to (a). The second term of the r.h.s. also approaches 0 as $\lambda \downarrow 0$ since $Q_\lambda \rightarrow Q$ as $\lambda \downarrow 0$.

(c) Note that, by definition of the perturbed t.p.f. P_λ , we have

$$P_\lambda R_\lambda = (1 - \varphi(\lambda)) P R_\lambda + \varphi(\lambda) Q_\lambda R_\lambda.$$

Note further that $Q_\lambda R_\lambda = P_\lambda^L$ and

$$(1 - \varphi(\lambda)) P R_\lambda = R_\lambda - \varphi(\lambda) I,$$

where I corresponds to the identity operator. Thus, we have

$$P_\lambda R_\lambda = R_\lambda - \varphi(\lambda) I + \varphi(\lambda) P_\lambda^L.$$

For any invariant probability measure of P_λ , μ_λ , we have

$$\mu_\lambda P_\lambda R_\lambda = \mu_\lambda R_\lambda - \varphi(\lambda)\mu_\lambda + \varphi(\lambda)\mu_\lambda P_\lambda^L,$$

which equivalently implies that

$$\mu_\lambda = \mu_\lambda P_\lambda^L,$$

since $\mu_\lambda P_\lambda = \mu_\lambda$. Thus, we conclude that μ_λ is also an invariant p.m. of P_λ^L .

(d) Let $\hat{\mu}$ denote a weak limit point of μ_λ as $\lambda \downarrow 0$. To see that such a limit exists, take $\hat{\mu}$ to be an invariant probability measure of P . Then,

$$\begin{aligned} & \|P_\lambda f - P f\|_\infty \\ & \geq \|\mu_\lambda(P_\lambda f - P f)\|_\infty \\ & = \|(\mu_\lambda - \hat{\mu})(I - P)[f]\|_\infty. \end{aligned}$$

Note that the weak convergence of P_λ to P , it necessarily implies that $\mu_\lambda \Rightarrow \hat{\mu}$. Note further that

$$\begin{aligned} & \hat{\mu}[f] - \hat{\mu}Q\Pi f \\ & = (\hat{\mu}[f] - \mu_\lambda[f]) + \mu_\lambda[P_\lambda^L f - Q\Pi f] + \\ & \quad (\mu_\lambda[Q\Pi f] - \hat{\mu}[Q\Pi f]). \end{aligned}$$

The first and the third term of the r.h.s. approaches 0 as $\lambda \downarrow 0$ due to the fact that $\mu_\lambda \Rightarrow \hat{\mu}$. The same holds for the second term of the r.h.s. due to part (b). Thus, we conclude that any weak limit point of μ_λ as $\lambda \downarrow 0$ is an invariant p.m. of $Q\Pi$. •

4.3 Invariant p.m. of one-step perturbed process

Define the finite-state Markov process \hat{P} as in (3).

Proposition 5 (Unique invariant p.m. of $Q\Pi$). *There exists a unique invariant probability measure $\hat{\mu}$ of $Q\Pi$. It satisfies*

$$\hat{\mu}(\cdot) = \sum_{s \in \mathcal{S}} \pi_s \delta_s(\cdot) \tag{9}$$

for some constants $\pi_s \geq 0$, $s \in \mathcal{S}$. Moreover, $\pi = (\pi_1, \dots, \pi_{|\mathcal{S}|})$ is an invariant distribution of \hat{P} , i.e., $\pi = \pi \hat{P}$.

Proof. From Proposition 3(d), we know that the support of Π is on the set of pure strategy states \mathcal{S} . Thus, the support of $Q\Pi$ is also on \mathcal{S} . From Proposition 4, we know that $Q\Pi$ admits an invariant measure, say $\hat{\mu}$, whose support is also \mathcal{S} . Thus, $\hat{\mu}$ admits the form of (9), for some constants $\pi_s \geq 0$, $s \in \mathcal{S}$.

Note also that $\mathcal{N}_\delta(s')$ is a continuity set of $Q\Pi(s, \cdot)$, i.e., $Q\Pi(s, \partial\mathcal{N}_\delta(s')) = 0$. Thus, from Portmanteau theorem, given that $QP^t \Rightarrow Q\Pi$,

$$Q\Pi(s, \mathcal{N}_\delta(s')) = \lim_{t \rightarrow \infty} QP^t(s, \mathcal{N}_\delta(s')) = \hat{P}_{ss'}.$$

If we also define $\pi_s \doteq \hat{\mu}(\mathcal{N}_\delta(s))$, then

$$\pi_{s'} = \hat{\mu}(\mathcal{N}_\delta(s')) = \sum_{s \in \mathcal{S}} \pi_s Q\Pi(s, \mathcal{N}_\delta(s')) = \sum_{s \in \mathcal{S}} \pi_s \hat{P}_{ss'},$$

which shows that π is an invariant distribution of \hat{P} , i.e., $\pi = \pi\hat{P}$.

It remains to establish uniqueness of the invariant distribution of $Q\Pi$. Note that the set \mathcal{S} of pure strategy states is isomorphic with the set \mathcal{A} of action profiles. If agent i trembles (as t.p.f. Q dictates), then all actions in \mathcal{A}_i have positive probability of being selected, i.e., $Q(\alpha, (\alpha'_i, \alpha_{-i})) > 0$ for all $\alpha'_i \in \mathcal{A}_i$ and $i \in \mathcal{I}$. It follows by Proposition 2 that $Q\Pi(\alpha, (\alpha'_i, \alpha_{-i})) > 0$ for all $\alpha'_i \in \mathcal{A}_i$ and $i \in \mathcal{I}$. Finite induction then shows that $(Q\Pi)^n(\alpha, \alpha') > 0$ for all $\alpha, \alpha' \in \mathcal{A}$. It follows that if we restrict the domain of $Q\Pi$ to \mathcal{S} , it defines an irreducible stochastic matrix. Therefore, $Q\Pi$ has a unique invariant distribution. •

4.4 Proof of Theorem 1

Theorem 1(a)–(b) is a direct implication of Propositions 4–5.

5 Conclusions & Future Work

In this paper, we considered a class of reinforcement-learning algorithms that belong to the family of learning automata, and we provided an explicit characterization of the invariant probability measure of its induced Markov chain. Through this analysis, we demonstrated convergence (in a weak sense) to the set of pure-strategy states, overcoming prior restrictions necessary under an ODE-approximation analysis, such as the existence of a potential function. Thus, we opened up new possibilities for equilibrium selection through this type of algorithms that goes beyond the fine class of potential games.

Although the set of pure-strategy-states (which are the stochastically-stable states) may contain non-Nash pure strategy profiles, a follow-up analysis that excludes convergence to such pure-strategy-states may be performed (similarly to the analysis presented in [11] for diminishing step size).

References

1. M. Tsetlin, *Automaton Theory and Modeling of Biological Systems*. Academic Press, 1973.

2. K. Narendra and M. Thathachar, *Learning Automata: An introduction*. Prentice-Hall, 1989.
3. W. B. Arthur, "On designing economic agents that behave like human agents," *J. Evolutionary Econ.*, vol. 3, pp. 1–22, 1993.
4. T. Börgers and R. Sarin, "Learning through reinforcement and replicator dynamics," *J. Econ. Theory*, vol. 77, no. 1, pp. 1–14, 1997.
5. I. Erev and A. Roth, "Predicting how people play games: reinforcement learning in experimental games with unique, mixed strategy equilibria," *Amer. Econ. Rev.*, vol. 88, pp. 848–881, 1998.
6. E. Hopkins and M. Posch, "Attainability of boundary points under reinforcement learning," *Games Econ. Behav.*, vol. 53, pp. 110–125, 2005.
7. A. Beggs, "On the convergence of reinforcement learning," *J. Econ. Theory*, vol. 122, pp. 1–36, 2005.
8. M. Thathachar and P. Sastry, *Networks of Learning Automata: Techniques for Online Stochastic Optimization*. Kluwer Academic Publishers, 2004.
9. G. Chasparis and J. Shamma, "Distributed dynamic reinforcement of efficient outcomes in multiagent coordination and network formation," *Dynamic Games and Applications*, vol. 2, no. 1, pp. 18–50, 2012.
10. D. Leslie, "Reinforcement learning in games," Ph.D. dissertation, School of Mathematics, University of Bristol, 2004.
11. G. C. Chasparis, J. S. Shamma, and A. Rantzer, "Nonconvergence to saddle boundary points under perturbed reinforcement learning," *Int. J. Game Theory*, vol. 44, no. 3, pp. 667–699, 2015.
12. M. Posch, "Cycling in a stochastic learning algorithm for normal form games," *J. Evolutionary Econ.*, vol. 7, pp. 193–207, 1997.
13. P. Sastry, V. Phansalkar, and M. Thathachar, "Decentralized learning of Nash equilibria in multi-person stochastic games with incomplete information," *IEEE Trans. Syst. Man Cybern.*, vol. 24, no. 5, pp. 769–777, 1994.
14. K. Verbeeck, A. Now, J. Parent, and K. Tuyls, "Exploring selfish reinforcement learning in repeated games with stochastic rewards," *Autonomous Agents and Multi-Agent Systems*, vol. 14, no. 3, pp. 239–269, Apr. 2007.
15. J. Hu and M. P. Wellman, "Nash Q-learning for general-sum stochastic games," *J. Machine Learning Research*, vol. 4, no. Nov, pp. 1039–1069, 2003.
16. M. L. Littman, "Markov games as a framework for multi-agent reinforcement learning," in *Proc. Int. Conf. Machine Learning*. Morgan Kaufmann, 1994, pp. 157–163.
17. A. C. Chapman, D. S. Leslie, A. Rogers, and N. R. Jennings, "Convergent Learning Algorithms for Unknown Reward Games," *SIAM J. Control Optim.*, vol. 51, no. 4, pp. 3154–3180, Jan. 2013.
18. D. Monderer and L. Shapley, "Potential games," *Games Econ. Behav.*, vol. 14, pp. 124–143, 1996.
19. D. Leslie and E. Collins, "Individual Q-Learning in Normal Form Games," *SIAM J. Control Optim.*, vol. 44, no. 2, pp. 495–514, Jan. 2005.
20. J. R. Marden, H. P. Young, G. Arslan, and J. S. Shamma, "Payoff based dynamics for multi-player weakly acyclic games," *SIAM J. Control Optim.*, vol. 48, no. 1, pp. 373–396, 2009.

21. G. Chasparis, A. Arapostathis, and J. Shamma, “Aspiration learning in coordination games,” *SIAM J. Control and Optim.*, vol. 51, no. 1, 2013.
22. O. Hernandez-Lerma and J. B. Lasserre, *Markov Chains and Invariant Probabilities*. Birkhauser Verlag, 2003.